

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2018-139043  
(P2018-139043A)

(43) 公開日 平成30年9月6日(2018.9.6)

(5) Int.Cl.		F I	テーマコード (参考)		
<b>G06F</b>	<b>19/24</b>	<b>(2011.01)</b>	G06F	19/24	4B063
C12N	15/09	(2006.01)	C12N	15/00	A
C12Q	1/68	(2018.01)	C12Q	1/68	Z

審査請求 未請求 請求項の数 8 O L (全 17 頁)

(21) 出願番号 特願2017-33381 (P2017-33381)  
(22) 出願日 平成29年2月24日 (2017. 2. 24)

(出願人による申告) (1) 平成27年度、国立研究開発法人科学技術振興機構、戦略的創造研究推進事業「野外環境と超並列高度制御環境の統合モデリングによる頑健性限界の解明と応用」委託研究、産業技術力強化法第19条の適用を受ける特許出願 (2) 平成27年度、国立研究開発法人科学技術振興機構、戦略的創造研究推進事業「共生ネットワークの分子基盤とその応用展開」委託研究、産業技術力強化法第19条の適用を受ける特許出願

(71) 出願人 597065329  
学校法人 龍谷大学  
京都府京都市伏見区深草塚本町67番地  
(74) 代理人 110000914  
特許業務法人 安富国際特許事務所  
(72) 発明者 岩山 幸治  
滋賀県大津市瀬田大江町横谷1-5 学校法人 龍谷大学内  
(72) 発明者 永野 惇  
滋賀県大津市瀬田大江町横谷1-5 学校法人 龍谷大学内  
Fターム(参考) 4B063 QA01 QA13 QQ42 QS39

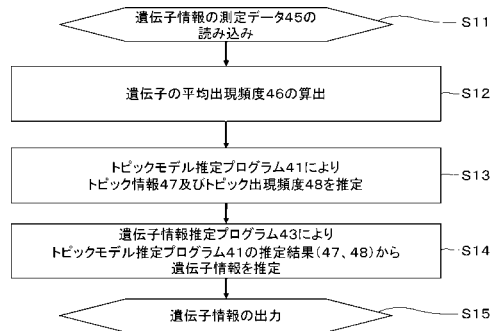
(54) 【発明の名称】 トピックモデルを用いた遺伝子情報推定装置

(57) 【要約】 (修正有)

【課題】低コストかつ高精度で遺伝子情報を推定できる装置を提供する。

【解決手段】既存のトピックモデルに基づき、遺伝子情報に適した新しいトピックモデルを用いた遺伝子情報を推定する。このモデルでは単語の生成分布を従来の多項分布から負の二項分布に置き換える。

【選択図】図2



## 【特許請求の範囲】

## 【請求項 1】

( i ) 複数のサンプルに対応する遺伝子情報の測定データを読み込む手順、  
 ( i i ) 遺伝子情報の測定データに基づき遺伝子の平均出現頻度を算出する手順、  
 ( i i i ) 各サンプルの各遺伝子にトピックが割り当てられる確率を初期化する手順、  
 ( i v ) 遺伝子へのトピック割り当て確率を更新する手順、  
 ( v ) トピック情報を更新する手順、  
 ( v i ) トピック情報の分散を更新する手順、  
 ( v i i ) 負の二項分布の  $d i s p e r s i o n$  パラメータを更新する手順、  
 および、  
 ( v i i i ) トピック情報およびトピック出現頻度から遺伝子情報を推定する手順、  
 を有する、  
 トピックモデルを用いた遺伝子情報推定装置。

10

## 【請求項 2】

手順 ( i v )、( v )、および ( v i i ) において、過分散の程度を表すパラメータを  $\phi_i$  とし、サンプル ( s ) においてトピック ( k ) を割り当てられた遺伝子 ( i ) の期待出現数を  $\mu_{s k i}$  としたとき、サンプル ( s ) における遺伝子 ( i ) の出現数 (  $r_{s i}$  ) が

## 【数 1】

$$P(r_{si}) = \frac{\Gamma(r_{si} + \phi_i^{-1})}{r_{si}! \Gamma(\phi_i^{-1})} \left( \frac{\mu_{ski}}{\mu_{ski} + \phi_i^{-1}} \right)^{r_{si}} \left( \frac{1}{1 + \mu_{ski} \phi_i} \right)^{\phi_i^{-1}}$$

20

という確率に従う、請求項 1 に記載の装置。

## 【請求項 3】

サンプル ( s ) における全遺伝子のカウントの総和 (  $\sum_s$  )、全サンプルでの全遺伝子の出現数の総和に対する遺伝子 ( i ) の出現数の割合の対数 (  $m_i$  )、指数分布あるいはジェフリーズ事前分布に従う分散 (  $\sigma_{k i}$  )、および平均が 0 で分散が  $\sigma_{k i}$  の正規分布に従うトピック ( k ) における遺伝子 ( i ) の期待カウントの平均からのずれを表す量を与えられたときに、サンプル ( s ) においてトピック ( k ) を割り当てられた遺伝子 ( i ) の期待出現数 (  $\mu_{s k i}$  ) が

30

$$\mu_{s k i} = \sum_s \exp( \sigma_{k i} + m_i )$$

と表される、請求項 2 に記載の装置。

## 【請求項 4】

請求項 1 ~ 3 のいずれか一項に記載の遺伝子情報推定装置を構成する各手段としてコンピュータを機能させるための、遺伝子情報推定プログラム。

## 【請求項 5】

請求項 4 に記載の遺伝子情報推定プログラムが記録されたことを特徴とするコンピュータ読み取り可能な記録媒体。

## 【請求項 6】

( i ) 複数のサンプルに対応する遺伝子情報の測定データを読み込むステップ、  
 ( i i ) 遺伝子情報の測定データに基づき遺伝子の平均出現頻度を算出するステップ、  
 ( i i i ) 各サンプルの各遺伝子にトピックが割り当てられる確率を初期化するステップ、  
 ( i v ) 遺伝子へのトピック割り当て確率を更新するステップ、  
 ( v ) トピック情報を更新するステップ、  
 ( v i ) トピック情報の分散を更新するステップ、  
 ( v i i ) 負の二項分布の  $d i s p e r s i o n$  パラメータを更新するステップ、  
 および、  
 ( v i i i ) トピック情報およびトピック出現頻度から遺伝子情報を推定するステップ、  
 を有する、トピックモデルを用いた遺伝子情報推定方法。

40

50

## 【請求項 7】

ステップ ( i v )、( v )、および ( v i i )において、過分散の程度を表すパラメータを  $\phi_i$  とし、サンプル ( s )においてトピック ( k )を割り当てられた遺伝子 ( i )の期待出現数を  $\mu_{s k i}$ としたとき、サンプル ( s )における遺伝子 ( i )の出現数 (  $r_{s i}$  )が、

## 【数 2】

$$P(r_{si}) = \frac{\Gamma(r_{si} + \phi_i^{-1})}{r_{si}! \Gamma(\phi_i^{-1})} \left( \frac{\mu_{ski}}{\mu_{ski} + \phi_i^{-1}} \right)^{r_{si}} \left( \frac{1}{1 + \mu_{ski} \phi_i} \right)^{\phi_i^{-1}}$$

という確率に従う、請求項 6 に記載の方法。

10

## 【請求項 8】

サンプル ( s )における全遺伝子のカウントの総和 (  $n_s$  )、全サンプルでの全遺伝子の出現数の総和に対する遺伝子 ( i )の出現数の割合の対数 (  $m_i$  )、指数分布あるいはジェフリーズ事前分布に従う分散 (  $\sigma_{k i}$  )、および平均が 0 で分散が  $\sigma_{k i}$  の正規分布に従うトピック ( k )における遺伝子 ( i )の期待カウントの平均からのずれを表す量が与えられたときに、サンプル ( s )においてトピック ( k )を割り当てられた遺伝子 ( i )の期待出現数 (  $\mu_{s k i}$  )が

$$\mu_{s k i} = n_s \exp(-\sigma_{k i}^2 / 2) \exp(\sigma_{k i}^2 / 2 + m_i)$$

と表される、請求項 7 に記載の方法。

20

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、トピックモデルを用いた遺伝子情報推定装置に関する。

## 【背景技術】

## 【0002】

近年、核酸塩基配列を次世代シーケンサで決定することで、遺伝子情報を網羅的に解析する技術が開発されている。例えば、DNA 上の遺伝子は、転写、翻訳などを経てその機能を発現するが、この遺伝子発現の状態を解析する目的で、転写産物である RNA の配列を次世代シーケンサで決定して遺伝子の発現を網羅的に定量する RNA - S e q 法が用いら

30

## 【0003】

RNA - S e q 等の遺伝子情報には、数千から数万の遺伝子の情報が含まれるため、出力されるデータは高次元のカウントデータとなる。同様の高次元カウントデータを扱うことの多い自然言語処理の分野では、文書の潜在的な意味を扱う手法としてトピックモデルが提案されている。トピックモデルは、文書内の単語の共起に基づき、文書を構成するトピックと各トピックにおける単語の出現頻度を同時に推定する。トピックモデルを RNA - S e q で得られた遺伝子発現データに適用することで、類似した発現パターンを持つ遺伝子群をトピックとして抽出できるのではないかと考えられる。実際、トピックモデルの一つである Latent Dirichle Allocation ( LDA ; 非特許文献 2 ) を、単一細胞の RNA - S e q データに適用することで、推定されたトピックから細胞間の階層的な構造を推定できることが示されている ( 非特許文献 3 )。

40

## 【先行技術文献】

## 【非特許文献】

## 【0004】

【非特許文献 1】 Wang et al. Nature Reviews Genetics 2009 Jan; 10 ( 1 ) : 57

【非特許文献 2】 Blei et al. Journal of Machine Learning Research 2003 3 : 993

【非特許文献 3】 DuVerle et al. BMC Bioinformatics 50

50

2016 Sep 13; 17(1): 363

【発明の概要】

【発明が解決しようとする課題】

【0005】

RNA-Seq等の遺伝子情報においては、読んだ配列の数（総リード数；depth）が多いほどノイズの少ない質の高いデータが得られる反面、コストが上がるために扱えるサンプルの数が減るというトレードオフが存在する。総リード数の少ないデータから精度よく発現量を推定できれば、より多くのサンプルの遺伝子情報発現量を定量することができる。RNA-Seqの場合、実際の発現データを見ると、発現パターンの関連した遺伝子が多数存在する。こうした遺伝子間の発現パターンの関係性をトピックモデルにより抽出することで、ノイズの大きな定量データからでも他の遺伝子の情報を利用して高精度な発現量の推定が可能になると期待される。しかし、RNA-Seqデータには、過分散、つまりポアソン分布や二項分布で予想される分散よりも大きな分散を持つという性質がある。従来のトピックモデルでは、多項分布を用いているため、このような大きな分散を説明することができない。

10

【課題を解決するための手段】

【0006】

本発明では、既存のトピックモデルに基づき、過分散を有する遺伝子情報の測定データにも適用可能な新しいトピックモデル及びその推定方法を提案する。このモデルでは、従来、単語の生成分布に用いられていた多項分布を、負の二項分布に置き換える。

20

【0007】

すなわち、本発明は、

- (i) 複数のサンプルに対応する遺伝子情報の測定データを読み込む手順、
- (ii) 遺伝子情報の測定データに基づき遺伝子の平均出現頻度を算出する手順、
- (iii) 各サンプルの各遺伝子にトピックが割り当てられる確率を初期化する手順、
- (iv) 遺伝子へのトピック割り当て確率を更新する手順、
- (v) トピック情報を更新する手順、
- (vi) トピック情報の分散を更新する手順、
- (vii) 負の二項分布のdispersionパラメータを更新する手順、

および、

- (viii) トピック情報およびトピック出現頻度から遺伝子情報を推定する手順、

30

を有する、トピックモデルを用いた遺伝子情報推定装置に関する。

【0008】

手順(iv)、(v)、および(vii)において、過分散の程度を表すパラメータを $\phi_i$ とし、サンプル(s)においてトピック(k)を割り当てられた遺伝子(i)の期待出現数を $\mu_{ski}$ としたとき、サンプル(s)における遺伝子(i)の出現数( $r_{si}$ )が

【数1】

$$P(r_{si}) = \frac{\Gamma(r_{si} + \phi_i^{-1})}{r_{si}! \Gamma(\phi_i^{-1})} \left( \frac{\mu_{ski}}{\mu_{ski} + \phi_i^{-1}} \right)^{r_{si}} \left( \frac{1}{1 + \mu_{ski} \phi_i} \right)^{\phi_i^{-1}}$$

40

という確率に従うことが好ましい。

【0009】

サンプル(s)における全遺伝子のカウントの総和( $\sum_i r_{si}$ )、全サンプルでの全遺伝子の出現数の総和に対する遺伝子(i)の出現数の割合の対数( $m_i$ )、指数分布あるいはジェフリーズ事前分布に従う分散( $\sigma_{ki}$ )、および平均が0で分散が $\sigma_{ki}$ の正規分布に従うトピック(k)における遺伝子(i)の期待カウントの平均からのずれを表す量を与えられたときに、サンプル(s)においてトピック(k)を割り当てられた遺伝子(i)の期待出現数( $\mu_{ski}$ )が

50

$\mu_{s k i} = s \exp (k_i + m_i)$   
と表される。

【0010】

また、本発明は、前記遺伝子情報推定装置を構成する各手段としてコンピュータを機能させるための、遺伝子情報推定プログラムに関する。

【0011】

また、本発明は、前記遺伝子情報推定プログラムが記録されたことを特徴とするコンピュータ読み取り可能な記録媒体に関する。

【0012】

また、本発明は、

(i) 複数のサンプルに対応する遺伝子情報の測定データを読み込むステップ、  
(ii) 遺伝子情報の測定データに基づき遺伝子の平均出現頻度を算出するステップ、  
(iii) 各サンプルの各遺伝子にトピックが割り当てられる確率を初期化するステップ、

(iv) 遺伝子へのトピック割り当て確率を更新するステップ、

(v) トピック情報を更新するステップ、

(vi) トピック情報の分散を更新するステップ、

(vii) 負の二項分布の dispersion パラメータを更新するステップ、

および、

(viii) トピック情報およびトピック出現頻度から遺伝子情報を推定するステップ、  
を有する、トピックモデルを用いた遺伝子情報推定方法に関する。

【0013】

ステップ (iv)、(v)、および (vii) において、過分散の程度を表すパラメータを  $\phi_i$  とし、サンプル (s) においてトピック (k) を割り当てられた遺伝子 (i) の期待出現数を  $\mu_{s k i}$  としたとき、サンプル (s) における遺伝子 (i) の出現数 ( $r_{s i}$ ) が、

【数2】

$$P(r_{si}) = \frac{\Gamma(r_{si} + \phi_i^{-1})}{r_{si}! \Gamma(\phi_i^{-1})} \left( \frac{\mu_{ski}}{\mu_{ski} + \phi_i^{-1}} \right)^{r_{si}} \left( \frac{1}{1 + \mu_{ski} \phi_i} \right)^{\phi_i^{-1}}$$

という確率に従うことが好ましい。

【0014】

サンプル (s) における全遺伝子のカウントの総和 ( $s$ )、全サンプルでの全遺伝子の出現数の総和に対する遺伝子 (i) の出現数の割合の対数 ( $m_i$ )、指数分布あるいはジェフリーズ事前分布に従う分散 ( $k_i$ )、および平均が 0 で分散が  $k_i$  の正規分布に従うトピック (k) における遺伝子 (i) の期待カウントの平均からのずれを表す量を与えられたときに、サンプル (s) においてトピック (k) を割り当てられた遺伝子 (i) の期待出現数 ( $\mu_{s k i}$ ) が

$$\mu_{s k i} = s \exp (k_i + m_i)$$

と表される。

【発明の効果】

【0015】

過分散のある遺伝子情報を扱うための新しいトピックモデルを提案し、シミュレーションによって推定精度の検証を行った。低コストを想定した総リード数が  $10^6$  程度のデータから本発明のモデルによって推定した遺伝子発現量は、よりコストのかかる総リード数  $10^9$  程度で定量した発現量よりも真の値に近かった (図 8)。このことは、1 サンプル当たりの総リード数を増やすことで質の良いデータを得るよりも、総リード数を減らして多くのサンプルの定量を行ったうえで本発明のモデルによる発現量の推定を行うことが有効であることを示している。

【図面の簡単な説明】

10

20

30

40

50

【0016】

【図1】本発明の実施形態に係る遺伝子情報推定装置1の構成を示すブロック図を示す。

【図2】遺伝子発現推定装置1による処理の流れを示すフローチャートを示す。

【図3A】トピックモデル推定プログラム41の一実施形態による処理の流れを示すフローチャートを示す。

【図3B】トピックモデル推定プログラム41の他の実施形態による処理の流れを示すフローチャートを示す。

【図4】遺伝子情報推定プログラム43による処理の流れを示すフローチャートを示す。

【図5】実施例における対数周辺尤度の下限の近似値を示す。

【図6】実施例において推定された各サンプルにおけるトピックの出現確率を示す。

10

【図7】実施例における各遺伝子の出現確率のトピック間の差を示す。

【図8】実施例におけるリード数の期待値と推定値との誤差の分布を示す。

【発明を実施するための形態】

【0017】

[1] 遺伝子情報推定装置

本発明は、

- (i) 複数のサンプルに対応する遺伝子情報の測定データを読み込む手順、
- (ii) 遺伝子情報の測定データに基づき遺伝子の平均出現頻度を算出する手順、
- (iii) 各サンプルの各遺伝子にトピックが割り当てられる確率を初期化する手順、
- (iv) 遺伝子へのトピック割り当て確率を更新する手順、
- (v) トピック情報を更新する手順、
- (vi) トピック情報の分散を更新する手順、
- (vii) 負の二項分布の dispersion パラメータを更新する手順、

20

および、

- (viii) トピック情報およびトピック出現頻度から遺伝子情報を推定する手順、
- を有する、

トピックモデルを用いた遺伝子情報推定装置に関する。

【0018】

手順(i)では複数のサンプルに対応する遺伝子情報の測定データを読み込む。手順(ii)では遺伝子情報の測定データに基づき遺伝子の平均出現頻度を算出する。遺伝子の平均出現頻度の算出法については後述する。手順(iii)では各サンプルの各遺伝子にトピックが割り当てられる確率を初期化する。初期化方法については後述する。

30

【0019】

遺伝子へのトピック割り当て確率を更新する手順(iv)では、各サンプルにおいて各遺伝子のカウントがどのトピックから生成されたかについて事後分布の推定を行う。

【0020】

トピック情報を更新する手順(v)では、割り当てられたトピックに基づき、トピック情報、すなわち各トピックにおける遺伝子の出現頻度の平均からのずれを推定する。

【0021】

トピック情報の分散を更新する手順(vi)では、推定されたトピック情報に基づき、トピック情報の分散を更新する。

40

【0022】

負の二項分布の dispersion パラメータを更新する手順(vii)では、トピックごとに各遺伝子のカウントデータを生成する負の二項分布のパラメータを推定する。すなわち、手順(iv)~(v)で割り当てられたトピックと推定されたトピック情報から算出される期待値と測定データの差分に基づき、負の二項分布の dispersion パラメータを推定する。

【0023】

トピック情報およびトピック出現頻度から遺伝子情報を推定する手順(viii)では、トピックモデル推定プログラムにより推定されたトピック情報と平均出現頻度から各トピ

50

ックにおける遺伝子の出現頻度を算出し、同じくトピックモデル推定プログラムで推定された各サンプルにおけるトピックの出現頻度と掛け合わせて得られる各サンプルにおける遺伝子の出現頻度を遺伝子情報の推定として出力する。

【0024】

[2] 遺伝子情報の測定データ

遺伝子情報の測定データは、試料から核酸を調製する工程、核酸からライブラリーを調製し増幅する工程、および次世代シーケンサーによりライブラリーの配列決定を行う工程を経ることにより得られる。各工程における作業手順および条件は特に限定されず、通常の方法を用いることが可能である。試料は限定されず、植物細胞、微生物細胞、動物細胞等を使用可能である。核酸はDNAであってもよく、RNAであってもよい。遺伝子情報は、複数の遺伝子の発現や、複数の遺伝子の存在に関連する情報であれば特に限定されないが、本発明の遺伝子情報推定装置は、過分散を持つ遺伝子情報に好適に適用できる。

10

【0025】

遺伝子情報の測定データとして、RNA-Seqデータ、ChIP-Seqデータが挙げられる。RNA-Seqデータは、RNAを次世代シーケンサーで定性・定量して得られ、トランスクリプトーム解析に用いられる。ChIP-Seqデータは、クロマチン免疫沈降法と次世代シーケンサーによる配列決定を組み合わせ得られ、DNA結合タンパク質が結合するDNA配列を示す。ChIP-Seqデータの取得時に対象となるDNA結合タンパク質としては、ヒストン、転写因子、DNA修飾酵素が挙げられる。

20

【0026】

遺伝子情報の測定データとして、DNAのメチル化を示すデータも挙げられる。DNAのメチル化を示すデータは、パイサルファイト変換を行ったゲノムDNAを次世代シーケンサーで解析してメチル化シトシンを検出する方法や、抗メチルシトシン抗体やメチルCpGタンパク質によるDNAの沈降と次世代シーケンサーによる配列決定を組み合わせる方法により得ることができる。

【0027】

本発明の装置は、メタゲノム解析において遺伝子情報を推定するためにも使用可能である。メタゲノム解析は、複数の微生物を含む試料から核酸を抽出し、その核酸塩基配列をまとめて解読し、試料中に存在した微生物の構成を解析する手法である。

30

【0028】

遺伝子情報の測定データとして、RNA-Seqデータを用いた場合、出力される遺伝子情報は、リード数、すなわち遺伝子発現量の期待値である。これにより、試料におけるトランスクリプトームを高い精度で推定できる。

【0029】

通常、遺伝子情報の推定のためにはサンプルあたりのデータ量が多い必要があるが、本発明の装置によれば、RNA-Seqの場合にはサンプルあたり100,000から1000,000リード程度という少ないデータ量でも高い精度で遺伝子情報を推定できる。

【0030】

[3] 遺伝子情報推定装置の全体構成

図1は、本発明の実施形態に係る遺伝子情報推定装置1の構成を示すブロック図である。遺伝子情報推定装置1は、演算手段2と、入力手段3と、記憶手段4と、出力手段5とを備えている。各手段2~5はバスライン6に接続されている。

40

【0031】

演算手段2は、例えば、CPU(Central Processing Unit)およびRAM(Random Access Memory)から構成される主制御装置である。演算手段2は、トピックモデル推定部21と、遺伝子情報推定部22と、メモリ23を含む。

【0032】

入力手段3は、例えば、キーボード、マウス、ディスクドライブ装置などから構成される。記憶手段4は、例えば、一般的なハードディスク装置などから構成され、プログラム格

50

納部 40 と、データ格納部 44 とを含む。

【0033】

プログラム格納部 40 には、演算手段 2 で用いられるプログラムとして、トピックモデル推定プログラム 41、および遺伝子情報推定プログラム 43 とを記憶しておくことが可能である。演算手段 2 は、記憶手段 4 のプログラム格納部 40 から、トピックモデル推定プログラム 41、および遺伝子情報推定プログラム 43 をそれぞれ読み込み、メモリ 23 に格納し実行することで、トピックモデル推定部 21 と、遺伝子情報推定部 22 とをそれぞれ実現する。

【0034】

データ格納部 44 には、演算手段 2 で用いられる各種データとして、遺伝子情報の測定データ 45、遺伝子の平均出現頻度 46、トピック情報 47、トピック出現頻度 48 とを記憶する。遺伝子情報の測定データ 45 は入力手段 3 を介して入力され、記憶手段 4 のデータ格納部 44 に記憶される構成とすることが可能である。遺伝子の平均出現頻度 46、トピック情報 47、トピック出現頻度 48 は、演算手段 2 の演算処理結果を示すデータである。

10

【0035】

出力手段 5 は、例えば、グラフィックボード（出力インタフェース）およびそれに接続されたモニタである。モニタは、例えば、液晶ディスプレイ等から構成され、可視化を行った結果等を表示する。

【0036】

本発明の課題は、前述の遺伝子情報推定装置を構成する各手段としてコンピュータを機能させるための遺伝子情報推定プログラム、前記遺伝子情報推定プログラムが記録されたコンピュータ読み取り可能な記録媒体によっても解決できる。また、本発明の課題は、前述の遺伝子情報推定装置の手順と同様のステップ含む、トピックモデルを用いた遺伝子情報推定方法によっても解決できる。

20

【0037】

遺伝子情報推定装置は、一般的なコンピュータを、遺伝子情報推定プログラムにより動作させることで実現できる。このプログラムは、通信回線を介して提供することも可能であるし、CD-ROM等のコンピュータ読み取り可能な記録媒体に書き込んで配布することも可能である。コンピュータは、コンピュータ読み取り可能な記録媒体から遺伝子情報推定プログラムを読み出し、このプログラムに基づいた各機能を実現することができる。また、遺伝子情報推定プログラムをインストールされたコンピュータは、CPUが、ROM等に格納されたこのプログラムをRAMに展開することにより、遺伝子情報推定装置としての効果を奏することができる。

30

【0038】

[4] 遺伝子情報推定装置による処理の流れ

図 2 は、図 1 に示した遺伝子情報推定装置 1 による処理の流れを示すフローチャートである。S11 において、遺伝子情報推定装置 1 は遺伝子情報の測定データ 45 を読み込む。S12 において、遺伝子情報推定装置 1 は遺伝子情報の測定データ 45 から遺伝子の平均出現頻度 46 を算出する。

40

【0039】

遺伝子の平均出現頻度 46 は、例えば遺伝子の全サンプルを通しての出現回数を、その全ての遺伝子についての総和で割ることで算出できる。

【0040】

S13 において、遺伝子情報推定装置 1 はトピックモデル推定プログラム 41 によりトピック情報 47 およびトピック出現頻度 48 を推定する。S14 において、遺伝子情報推定装置 1 はトピック情報 47 およびトピック出現頻度 48 から遺伝子情報を推定する。S15 において、遺伝子情報推定装置 1 は推定された遺伝子情報を出力手段 5 から出力する。

【0041】

[5] トピックモデル推定プログラムによる処理の流れ

50



図3Aは、図1に示したトピックモデル推定プログラム41による処理の流れを示すフローチャートである。

【0042】

S201において、遺伝子情報推定装置1は、トピックモデル推定プログラム41の初期化を行う。S202において、遺伝子情報推定装置1は遺伝子へのトピック割り当て確率の更新を行う。S203において、遺伝子情報推定装置1はトピック情報( )の更新1を行う。S204において、遺伝子情報推定装置1はトピック情報の分散( )の更新を行う。S205において、遺伝子情報推定装置1は負の二項分布の dispersion パラメータ( )の更新を行う。S206において、遺伝子情報推定装置1は対数周辺尤度の下限が収束しているか否かを評価し、収束していない場合にはS202に戻る。対数周辺尤度の下限が収束している場合にはS211に進む。S211において、遺伝子情報推定装置1はトピック情報47およびトピック出現頻度48を出力手段5から出力する。

10

【0043】

以下、トピックモデル推定プログラムにおける処理の流れの詳細を、必要に応じて図3Aを参照しながら説明する。過分散なデータを単一のトピックから発生させるために、先行技術である Sparse Additive Generative model (SAGE; Eisenstein et al. Proc. Int. Conf. Mach. Learn. 2011)における単語あるいは遺伝子の出現確率を負の二項分布に置き換える。サンプルで各遺伝子に割り当てられるトピックの確率は、SAGEと同様に Dirichlet 分布から生成される。

20

【0044】

【数3】

$$\theta_s \sim \text{Dir}(\alpha), s = 1, \dots, S$$

【数4】

$$z_{si} \sim \text{Multi}(\theta_s), i = 1, \dots, G$$

【0045】

ここで、 $\theta_s$  はサンプル  $s$  において各トピックを割り当てる確率、 $z_{si}$  はサンプル  $s$  の  $i$  番目の遺伝子に割り当てられたトピック、 $s$  と  $G$  はそれぞれサンプルと遺伝子の数を表す。先行技術と同様に、すべてのトピックにおける平均的な遺伝子の出現頻度の対数を  $m$ 、各トピックにおける出現頻度がこの  $m$  からどの程度ずれているかを表す量を  $\tau_k$  とし、その事前分布は以下の正規分布とする。

30

【0046】

【数5】

$$\eta_{k,i} \sim \mathcal{N}(0, \tau_{k,i})$$

【0047】

分散  $\tau_{k,i}$  は指数分布

【数6】

$$P(\tau_{k,i}) \propto \exp(-\gamma \tau_{k,i})$$

あるいは Jeffrey's 事前分布

【数7】

$$P(\tau_{k,i}) \propto 1/\tau_{k,i}$$

を用いる。こうすることで、 $\tau_k$  は多くの要素が0となることが知られている。

サンプル  $s$  における遺伝子  $i$  にトピック  $k$  が割り当てられたとき、すなわち  $z_{si} = k$  のとき、そのリード数は、

40

## 【数 8】

$$P(r_{si}) = \frac{\Gamma(r_{si} + \phi_i^{-1})}{r_{si}! \Gamma(\phi_i^{-1})} \left( \frac{\mu_{ski}}{\mu_{ski} + \phi_i^{-1}} \right)^{r_{si}} \left( \frac{1}{1 + \mu_{ski} \phi_i} \right)^{\phi_i^{-1}}$$

という負の二項分布に従うものとする。ここで、 $\phi_i$  は負の二項分布において過分散の度合いを決める dispersion というパラメータである。この分布の期待出現数  $\mu_{s, k, i}$  は、サンプル  $s$  における全遺伝子のカウントの総和を  $r_{si}$  と表記すると、

## 【数 9】

$$\mu_{ski} = \nu_s \lambda_{k,i} = \nu_s \exp(m_i + \eta_{k,i})$$

10

となる。

## 【0048】

トピックモデル推定プログラムの初期化処理 (S201) では、各サンプルの各遺伝子にトピックの割り当てられる確率をランダムに初期化する。また、入力されたカウントデータから平均的な出現頻度の対数  $m_i$  を算出し、上記のトピックの割り当てられる確率と合わせて平均出現頻度からのずれ  $\eta_{k,i}$  の初期値を決める。

## 【0049】

遺伝子へのトピック割り当て確率の更新 (S202) は下記の手順により行う。すなわち、RNA-Seq データに基づき、 $z$ 、 $\eta$ 、 $\tau$  の事後分布の推定を行う。これらの事後分布を解析的に求めることはできないため、事後分布の近似  $q(z)$ 、 $q(\eta)$ 、 $q(\tau)$  を考え、対数周辺尤度の下限

20

## 【数 10】

$$F(q(z), q(\eta), q(\tau)) = \mathbb{E}_{q(z), q(\eta)} [\log P(\mathbf{r}, \mathbf{z} | \boldsymbol{\alpha}, \mathbf{m}, \boldsymbol{\eta}, \boldsymbol{\phi}) + \log P(\boldsymbol{\eta} | \boldsymbol{\tau}) + \log P(\boldsymbol{\tau})] \\ - \mathbb{E}_{q(z)} [\log q(z)] - \mathbb{E}_{q(\eta)} [\log q(\eta)] - \mathbb{E}_{q(\tau)} [\log q(\tau)]$$

を最大化する。ここで

## 【数 11】

$$\mathbb{E}_q[\cdot]$$

30

は確率分布  $q(\cdot)$  のもとでの期待値を意味する。同様に確率分布  $q(\cdot)$  のもとでの分散を

## 【数 12】

$$\mathbb{V}_q[\cdot]$$

、サンプル  $s$  の遺伝子  $i$  を除いた  $z$  を  $z \setminus s, i$ 、サンプル  $s$  において遺伝子  $i$  以外でトピック  $k$  を割り当てられた遺伝子の数を  $n_{s, k} \setminus s, i$  と表記し、

## 【数 13】

$$\hat{\lambda}_{k,i} = \exp(m_i + \mathbb{E}_{q(\eta)}[\eta_{k,i}])$$

とおく。

40

## 【0050】

トピックモデルの推定は、 $\mu_{s, k, i}$  の事前分布を指数分布として推定を行う。最初に遺伝子へのトピックの割り当て確率の更新 (S202) を行う。トピックの割り当て確率  $q(z)$  の更新式は、

【数 1 4】

$$\begin{aligned} & \log q(z_{si} = k) \\ & \approx \log \left( \mathbb{E}_{q(z^{s,i})} \left[ n_{s,k}^{s,i} \right] + \alpha_k \right) - \frac{\mathbb{V}_{q(z^{s,i})} \left[ n_{s,k}^{s,i} \right]}{2 \left( \mathbb{E}_{q(z^{s,i})} \left[ n_{s,k}^{s,i} \right] + \alpha_k \right)^2} \\ & \quad - r_{si} \log \left( 1 + 1/\nu_s \hat{\lambda}_{k,i} \phi_i \right) - \phi_i^{-1} \log \left( 1 + \nu_s \hat{\lambda}_{k,i} \phi_i \right) \\ & \quad - \frac{(r_{si} + \phi_i^{-1}) \mathbb{V}_{q(\eta)} [\eta_{k,i}]}{2 \left( 1 + 1/\nu_s \hat{\lambda}_{k,i} \phi_i \right) \left( 1 + \nu_s \hat{\lambda}_{k,i} \phi_i \right)} + \text{const.} \end{aligned} \tag{10}$$

となる。

【0051】

次に、トピック情報 ( ) の更新 1 (S203) を行う。q ( ) はラプラス近似を行う。関数

【数 1 5】

$$f(\eta_{k,i}) = \mathbb{E}_{q(z), q(\eta)} \left[ -r_{si} \log \left( 1 + \frac{1}{\nu_s \lambda_{k,i} \phi_i} \right) - \phi_i^{-1} \log (1 + \nu_s \lambda_{k,i} \phi_i) \right] - \frac{\mathbb{E}_{q(\tau)} [\tau_{k,i}^{-1}] \mathbb{E}_{q(\eta)} [\eta_{k,i}^2]}{2} \tag{20}$$

を定義し、

【数 1 6】

$$\eta_{k,i}^* = \arg \max_{\eta_{k,i}} f(\eta_{k,i})$$

と表記すると

【数 1 7】

$$q(\eta_{k,i}) \approx \mathcal{N} \left( \eta_{k,i}^*, -\frac{1}{f''(\eta_{k,i}^*)} \right)$$

となる。ここで、

【数 1 8】

$$\eta_{k,i}^*$$

の推定値はニュートン法による更新を 1 ステップにつき 1 回行う、

【数 1 9】

$$\Delta \eta_{k,i}^* = -\frac{f'(\eta_{k,i}^*)}{f''(\eta_{k,i}^*)}$$

【0052】

次に、トピック情報の分散 ( ) の更新 (S204) を行う。 の事後分布 q ( ) はガンマ分布

【数 2 0】

$$q(\tau_{k,i} | a_{k,i}, b_{k,i}) = \tau_{k,i}^{a_{k,i}-1} \frac{\exp(-\tau_{k,i}/b_{k,i})}{\Gamma(a_{k,i}) b_{k,i}^{a_{k,i}}}$$

とし、そのパラメータ  $a_{k,i}$  と  $b_{k,i}$  を

【数 2 1】

$$-\Delta a_{k,i} = \frac{(\frac{1}{2} - a_{k,i})\psi_1(a_{k,i}) + \frac{\mathbb{E}_{q(\eta)}[\eta_{k,i}^2]}{2(a_{k,i}-1)^2 b_{k,i}} - \gamma b_{k,i} + 1}{(\frac{1}{2} - a_{k,i})\psi_2(a_{k,i}) - \psi_1(a_{k,i}) - \frac{\mathbb{E}_{q(\eta)}[\eta_{k,i}^2]}{(a_{k,i}-1)^3 b_{k,i}}}$$

【数 2 2】

$$b_{k,i} = \frac{1 + \sqrt{1 + 8\gamma \mathbb{E}_{q(\eta)}[\eta^2] a_{k,i}/(a_{k,i} - 1)}}{4\gamma a_{k,i}}$$

で更新する。

10

【0053】

以上の推定値から算出される期待出現頻度と測定されたカウントデータ  $r_{s_i}$  との差分に基づき、負の二項分布の dispersion パラメータの更新を行う (S205)。以上の更新 (S202 ~ S205) を対数周辺尤度の下限が収束するまで繰り返す (S206)。対数周辺尤度の下限は各更新ステップで導出された近似に基づき算出される。

【0054】

より高精度の推定を行うために、トピックモデルの推定は複数段階行ってもよい。トピックモデルの推定を二段階行う場合のトピック推定プログラム 41 による処理の流れを図 3B に示す。

【0055】

20

図 3B において、S201 ~ S206 および S211 は図 3A と同じである。S206 において、遺伝子情報推定装置 1 は対数周辺尤度の下限が収束しているか否かを評価し、収束していない場合には S202 に戻る。対数周辺尤度の下限が収束している場合には S207 に進む。S207 において、遺伝子情報推定装置 1 は遺伝子へのトピック割り当て確率の更新を行う。S208 において、遺伝子情報推定装置 1 はトピック情報 ( ) の更新 2 を行う。S209 において、遺伝子情報推定装置 1 は、負の二項分布の dispersion パラメータ ( ) の更新を行う。S210 において、遺伝子情報推定装置 1 は対数周辺尤度の下限が収束しているか否かを評価し、収束していない場合には S207 に戻る。対数周辺尤度の下限が収束している場合には S211 に進む。S211 において、遺伝子情報推定装置 1 はトピック情報 47 およびトピック出現頻度 48 を出力手段 5 から出力する。

30

【0056】

トピックモデルの推定を二段階行う場合、まず、 $k_i$  の事前分布を指数分布として前述の方法で推定を行う。対数周辺尤度の下限が収束した後、よりスパースな  $k_i$  を得るために、 $k_i$  の事前分布を Jeffrey's 事前分布として遺伝子へのトピック割り当て確率 (S207)、トピック情報 (S208)、負の二項分布の dispersion パラメータ (S209) の更新を、対数周辺尤度の下限がもう一度収束するまで繰り返す。q (z) は一段階目と同じ式で更新できる。また、q ( ) はキャンセルされて消えるため、二段階目では更新を行わない。q ( ) の更新のみ、

【数 2 3】

40

$$f'(\eta_{k,i})$$

$$f''(\eta_{k,i})$$

中の

【数 2 4】

$$\mathbb{E}_{q(\tau)} \left[ \tau_{k,i}^{-1} \right]$$

が、

## 【数 2 5】

$$1/\eta_{k,i}$$

となるために更新式が異なる (S 2 0 8)。

## 【0 0 5 7】

遺伝子へのトピック割り当て確率の更新、トピック情報 ( ) の更新 2、および負の二項分布の dispersion パラメータの更新 (S 2 0 7 ~ S 2 0 9) を、対数周辺尤度  
10  
の下限が収束するまで繰り返す (S 2 1 0)。対数周辺尤度  
の下限の収束は、S 2 0 6  
と同じ方法により評価できる。対数周辺尤度  
の下限が収束した後、トピック情報 4 7 および  
トピック出現頻度 4 8 が出力される (S 2 1 1)。

## 【0 0 5 8】

S 2 0 2 ~ S 2 0 5、および S 2 0 7 ~ S 2 0 9 において、dispersion は負  
の二項分布の分散と期待値と dispersion の関係、

## 【数 2 6】

$$V[r_{s,i}] = E[r_{s,i}] + E[r_{s,i}]^2 \phi_i$$

から、リード数の期待値と実測の二乗誤差に基づいて推定する。

## 【0 0 5 9】

[ 5 ] 遺伝子情報推定プログラムによる処理の流れ

図 4 は、図 1 に示した遺伝子情報推定プログラム 4 3 による処理の流れを示すフローチャートである。S 3 1 において、遺伝子情報推定装置 1 は、遺伝子の平均出現頻度 4 6、ト  
20  
ピック情報 4 7、およびトピック出現頻度 4 8 を読み込む。遺伝子の平均出現頻度 4 6 は  
図 2 の S 1 2 で算出されたものであり、トピック情報 4 7、およびトピック出現頻度 4 8  
は図 3 A および図 3 B の S 2 1 1 で出力されたものである。S 3 2 において、遺伝子情報  
推定装置 1 は遺伝子情報を推定する。S 3 3 において、遺伝子情報推定装置 1 は、推定さ  
れた遺伝子情報を出力手段 5 から出力する。

## 【0 0 6 0】

S 3 2 における遺伝子情報の推定において、サンプル ( s ) における遺伝子 ( i ) の期待  
出現数

## 【数 2 7】

$$\hat{\mu}_{si}$$

は、

サンプル ( s ) における全遺伝子のカウントの総和 (  $\sum_s$  )、トピックモデル推定プログラ  
ム ( S 1 3 ) により推定された全サンプルでの全遺伝子の出現数の総和に対する遺伝子  
( i ) の出現数の割合の対数 (  $m_i$  )、および、以上で推定された各サンプルにおける各  
遺伝子へのトピックの割り当て確率、トピック ( k ) における遺伝子 ( i ) の期待カウ  
ントの平均からのずれを表す量 ( ) から、

## 【数 2 8】

$$\hat{\mu}_{si} = \sum_k q(z_{ki}) E_{q(\eta)} [\lambda_{k,i}] = \sum_k q(z_{s,k}) \exp \left( m_i + \eta_{k,i}^* + \frac{V_{q(\eta)} [\eta_{k,i}]}{2} \right)$$

と推定される。

## 【実施例】

## 【0 0 6 1】

(実施例 1) トピックの出現確率の推定

先行技術との比較を行うために、疑似的な RNA - Seq データを作成して解析を行った  
。統制群と実験群がそれぞれ 2 0 サンプルある状況を想定した。遺伝子の総数を 1 0 , 0  
0 0 とし、そのうち 5 0 0 遺伝子は実験群で発現量が増加し、別の 5 0 0 遺伝子は発現量  
が減少する。残りの 9 , 0 0 0 遺伝子は二群間で発現量が変わらないとする。統制群にお  
30  
40  
50

ける各遺伝子のリード数は以下の負の二項分布に従う。

【 0 0 6 2 】

【 数 2 9 】

$$\log \bar{\mu}_i \sim \mathcal{N}(0, 1)$$

【 数 3 0 】

$$\mu_i^c = \frac{\bar{\nu} \bar{\mu}_i}{\sum_i \bar{\mu}_i}$$

【 数 3 1 】

$$P(r_{si}) = \frac{\Gamma(r_{si} + \phi_i^{-1})}{r_{si}! \Gamma(\phi_i^{-1})} \left( \frac{\mu_i^c}{\mu_i^c + \phi_i^{-1}} \right)^{r_{si}} \left( \frac{1}{1 + \mu_i^c \phi_i} \right)^{\phi_i^{-1}}$$

10

【 0 0 6 3 】

【 数 3 2 】

$\bar{\nu}$   
は総リード数の期待値で、低コストで総リード数の少ない条件

【 数 3 3 】

$$\bar{\nu} = 10^6$$

20

と高コストで総リード数の多い条件

【 数 3 4 】

$$\bar{\nu} = 10^8$$

の二条件を想定した。dispersion<sub>i</sub> は、公開されている実測のデータ (Pickrell et al. Nature 2010) に基づき edgeR (Zhou et al. Nucleic Acids Res. 2014) で推定を行った。実験群の期待リード数  $\mu_i^e$  は以下のように決定した。

【 0 0 6 4 】

【 数 3 5 】

$$\mu_i^e = \mu_i^c (\gamma_i + 1.5)^{\rho_i}$$

【 数 3 6 】

$$P(\gamma_i) \propto \exp(-\gamma_i)$$

30

【 0 0 6 5 】

ここで、 $\rho_i$  は実験群で発現量が増加する遺伝子群については、 $\rho_i = +1$ 、減少する遺伝子群については  $\rho_i = -1$ 、残りの変化しない遺伝子群については、 $\rho_i = 0$  とした。生成した疑似データについて、LDA、SAGE、本発明の提案モデルの3つのモデルを適用した。

40

【 0 0 6 6 】

異なるトピック数で提案モデルの推定を行った時の、対数周辺尤度の下限の近似値を図5に示す。図5の横軸はトピック数を表し、縦軸は変分下限の近似値を表す。トピック数が2の時に最大となることから、以降の解析ではトピック数を2とした。疑似データは統制群と実験群の二群からなるため、この結果はトピック数を適切に選択できていることを示している。

【 0 0 6 7 】

3つのモデルで推定された各サンプルにおけるトピックの出現確率を図6に示す。図6において、横軸はサンプルを表し、縦軸はトピックの出現確率を表す。統制群と実験群を左

50

右に分けて表示している。すべてのモデルで、二群間でトピックの出現確率が変化しているが、LDAとSAGEに対して、本発明のモデル ( p r o p o s e d ) では二群間の違いが顕著になっている。

【0068】

(実施例2)トピックにおける遺伝子の出現確率の比較

各遺伝子について推定された二つのトピックにおける出現確率の差を比較した。LDAについては、各トピックにおける出現確率の差を、SAGEと本発明のモデルについては、

【数37】

$$\eta_{2,i} - \eta_{1,i}$$

を計算し、ヒストグラムを作成した(図7)。図7において、実験群で発現量が増加する500遺伝子 ( u p ) を白色、減少する500遺伝子 ( d o w n ) を黒色、それ以外の9,000遺伝子 ( c o n s t ) を灰色で表示している。LDAでは、ほとんどの遺伝子が0付近に集中しており、発現量の増加する遺伝子と減少する遺伝子とそれ以外の遺伝子を区別することが困難である。3種類の遺伝子は、SAGEよりも本発明のモデル ( p r o p o s e d ) の方がより明瞭に分離している。

【0069】

(実施例3)発現量の推定精度の比較

発現量の推定精度を比較するために、リード数の期待値  $\mu_{s_i}$  と推定値の誤差を算出した(図8)。図8ではリード数の期待値の大きさに応じて分割してプロットした。期待総リード数  $10^6$  で生成した疑似データ ( r a w )、コストをかけた状態を想定して期待総リード数を  $10^8$  とした疑似データ ( d e e p )、SAGEの推定 ( s a g e )、本発明のモデルの推定 ( p r o p o s e d ) を、それぞれ黒色の実線、灰色の実線、黒色の破線、灰色の破線でプロットした。

【0070】

総リード数の多い状況を想定したデータ ( d e e p ) は総リード数をそろえるために、生成したデータを100で割った値と期待値を比較した。LDAは、元のデータ ( r a w ) よりも誤差が大きくなっており、発現量の推定には使用できない。sageと本発明のモデルを比較すると、約60%の観測で本発明のモデルの方が期待値に近い推定を与えている。

【0071】

本発明のモデルは、今回のシミュレーションの二群比較のように統制のとれた実験のサンプルだけではなく、e d g e R や D E S e q といった負の二項分布を用いた既存の解析法が利用できない未知の構成のサンプルから定量したRNA-Seqデータに対しても、有効と考えられる。

【符号の説明】

【0072】

- 1 遺伝子情報推定装置
- 2 演算手段
- 2 1 トピックモデル推定部
- 2 2 遺伝子情報推定部
- 2 3 メモリ
- 3 入力手段
- 4 記憶手段
- 4 0 プログラム格納部
- 4 1 トピックモデル推定プログラム
- 4 3 遺伝子情報推定プログラム
- 4 4 データ格納部
- 4 5 遺伝子情報の測定データ
- 4 6 遺伝子の平均出現頻度

10

20

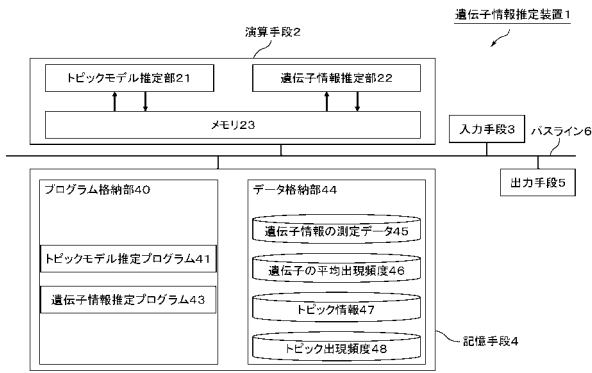
30

40

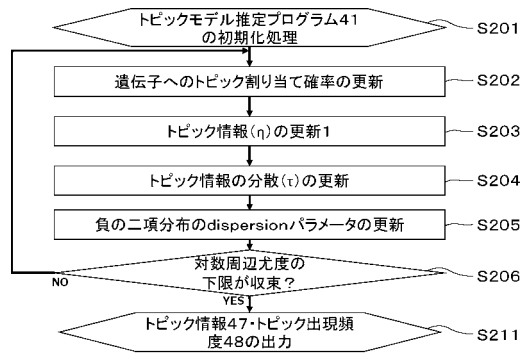
50

- 4 7 トピック情報
- 4 8 トピック出現頻度
- 5 出力手段
- 6 バスライン

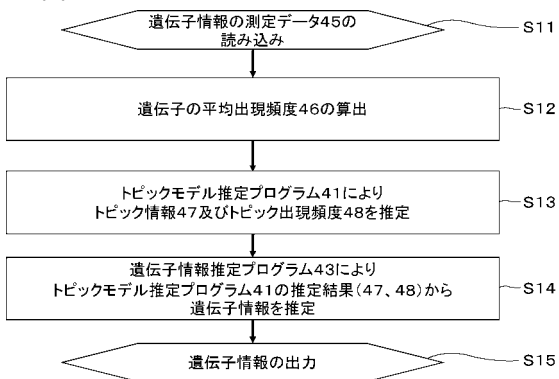
【 図 1 】



【 図 3 A 】

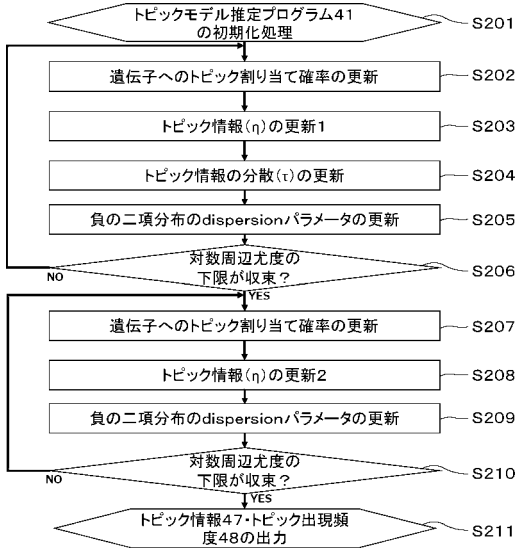


【 図 2 】

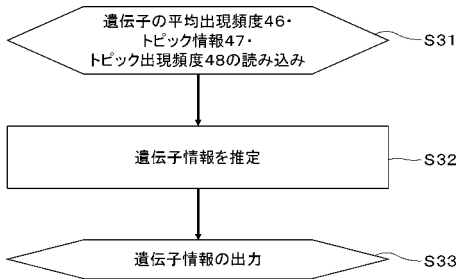




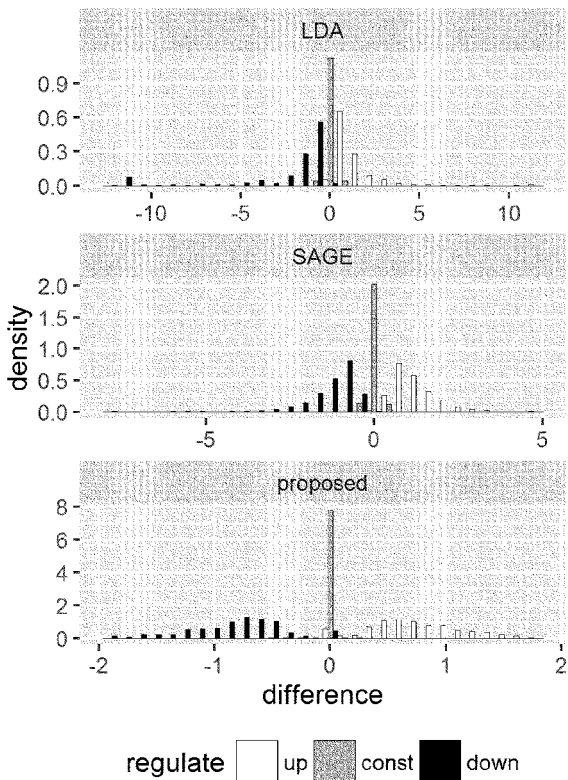
【 図 3 B 】



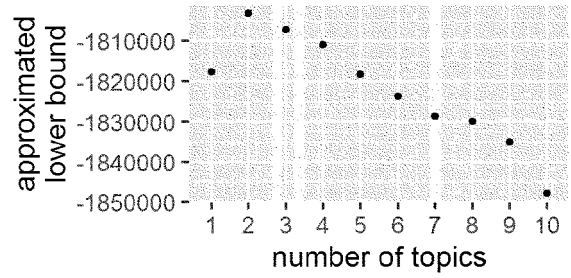
【 図 4 】



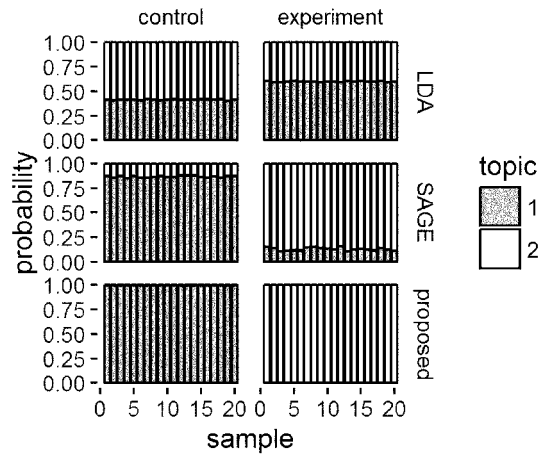
【 図 7 】



【 図 5 】



【 図 6 】



【 図 8 】

